

**Applications of Random Coding and Algebraic
Coding Theories to Universal Lossless Source
Coding Performance Bounds**

Gil I. Shamir

Department of Electrical & Computer Engineering

University of Utah

Salt Lake City, UT 84112

U.S.A.

DIMACS - 2003

Workshop on Algebraic Coding Theory and Information Theory
DIMACS Center, Rutgers University, Piscataway, NJ
December 15-18, 2003

Overview

Research Problem

- Average Case Universal Lossless Compression
- Performance Lower Bounds (on Redundancy - best possible performance of any scheme for a specific model)

Research Approach

- Use Redundancy-Capacity Theorems to obtain bounds
- Lower bound the relevant capacity for given source model

Models Discussed

- finite number of parameters parametric sources
- i.i.d. sources with large alphabets
- patterns induced by i.i.d. sources
- piecewise stationary sources
- piecewise stationary sources with slowly varying statistics
- switching sources

Universal Coding and Redundancy

Problem Layout

- A sequence x^n of length n , governed by P_θ ,
- θ unknown in a known class Λ ,
- uniquely decipherable code $L(\cdot)$ may depend on Λ but independent of θ .
- Unknown parameters cost redundancy.

Average Redundancy

of code $L(\cdot)$ for n -sequences drawn by source θ

$$R_n(L, \theta) \triangleq \frac{1}{n} E_\theta L(X^n) - H_\theta(X^n)$$

• E_θ - mean w.r.t. θ ,

• H_θ - per symbol entropy.

Average Universality Measure of a Class Λ

- **Maximin** $R_n^-(\Lambda)$ and **Minimax** $R_n^+(\Lambda)$ average redundancies - best code for some worst average (over x^n) case. [Davisson, 1973]
- **Average redundancy for most sources** [Rissanen, 1984] (strongest sense).

Redundancy-Capacity Theorem

Weak Version [Implied from Davission, 1973, Gallager, 1976]

Let $n \rightarrow \infty$. Let φ be a set of M points θ in the class Λ_k , that are *distinguishable* by x^n . Then, the minimax and maximin redundancies satisfy

$$R_+^n(\Lambda_k) = R_-^n(\Lambda_k) \geq (1 - \varepsilon) \frac{\log M}{n}$$

Strong Random Coding Version [Merhav & Feder, 1995, 1996]

Let $n \rightarrow \infty$. Define a distribution over Λ_k , and partition *most* of the class Λ^ε into disjoint countable sets φ , where the marginal of each $\theta \in \varphi$ is equal, and there are $M_\varphi \geq M$ sources in φ , *distinguishable* by x^n . Then,

$$R_n(L, \theta) \geq (1 - \varepsilon) \frac{\log M}{n},$$

for every code $L(\cdot)$, and almost every $\theta \in \Lambda_k$.

Distinguishability

θ and θ' distinguishable if x^n generated by θ appears to be generated by θ' with probability that goes to 0 and vice versa.

Use of Redundancy-Capacity Theorem

Weak Version for Λ_k

1. Demonstrate how to find φ .

2. Lower bound M .

3. Prove that all $\varphi \in \theta$ are distinguishable by x^n .

Strong Version for Λ_k

1. Demonstrate how to define most of the class Λ_ε .

2. Show that Λ_ε is most of the class.

3. Show how to partition Λ_ε such that every source in Λ_ε is in **exactly** one φ , and sources in φ are uniformly distributed with the uniform prior on Λ_k .

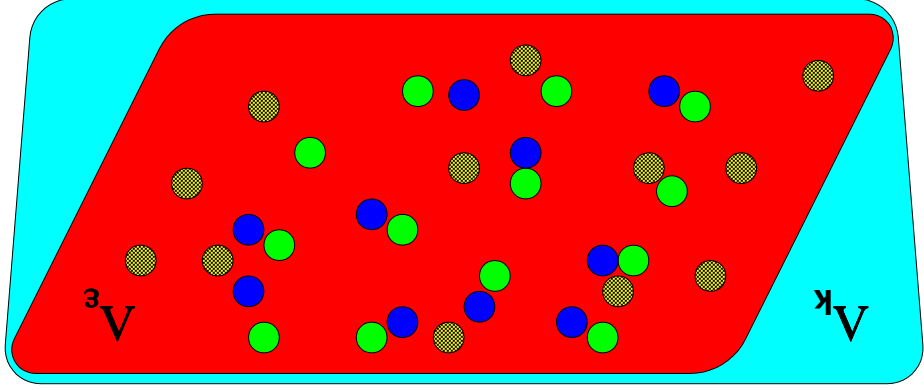
Lower bound M .

4. Prove that for every valid φ , all $\varphi \in \theta$ are distinguishable by x^n .

Compound Classes

If $\Lambda = \bigcup_k \Lambda_k$, redundancy for $\theta \in \Lambda_k$ consists of Intra-class redundancy in Λ_k , and Inter-class redundancy distinguishing Λ_k from Λ .

Redundancy Capacity - Demo



- $\theta \in \Phi_1$ $M^{\Phi_1} = 13$
 - $\theta \in \Phi_2$ $M^{\Phi_2} = 10$
 - $\theta \in \Phi_3$ $M^{\Phi_3} = 12$
- M = 10**

- The volume of Λ_k outside Λ_ϵ assumed negligible.
- Any θ is contained in a unique φ and has equal probability to other $\theta' \in \varphi$.
- In every φ all points distinguishable by x^n .

By theorem, for every code and almost every $\theta \in \Lambda_k$,

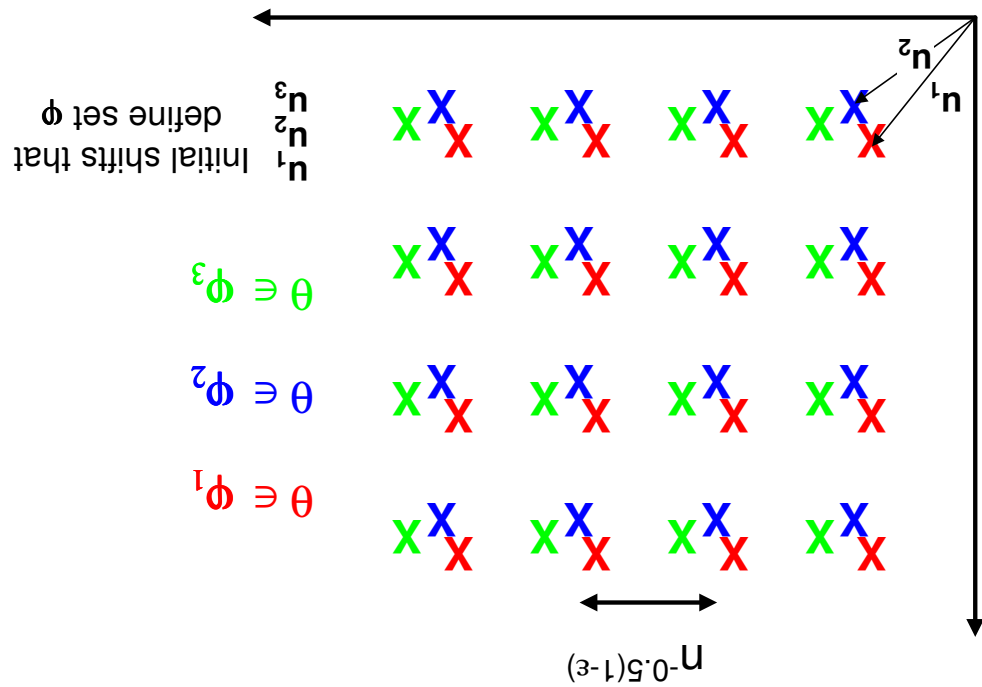
$$R_n(L, \theta) \geq (1 - \epsilon) \frac{n}{\log 10}$$

Finite k -dimensional Parametric Sources

- φ determined by initial shift \mathbf{u} in a grid (one φ sufficient for maximin)
- $\theta \in \varphi$ distinguishable if φ is a grid with spacing $n^{-0.5(1-\epsilon)}$

$$R_n(L, \theta) \geq (1 - \epsilon) \frac{k \log n}{2n}$$

for every code $L(\cdot)$ and almost every $\theta \in \Lambda$ [Rissanen, 1984]



Distinguishability

Setting and Proof in most sources sense

- Choose a random grid φ (as in random coding).
- Generate x^n by a given $\theta \in \varphi$.
- Let $\hat{\theta}$ be the Maximum Likelihood estimator of θ from x^n .
- Let $\hat{\theta}_g$ be the grid point whose components are nearest $\hat{\theta}$.
- Prove that $P_e = \Pr(\hat{\theta}_g \neq \theta \mid \theta) \rightarrow 0$ as $n \rightarrow \infty$.

Use union bound on components of θ :

$$P_e \leq \sum_k^{j=1} \Pr(\hat{\theta}^{g_i} \neq \theta_i)$$

$$\leq \sum_k n \cdot 2^{-n \cdot \min_{x_n \in A_i} D(P_{\theta_i} \| P_{\theta_i})}$$

$$\leq 2^{(\log k) + (\log n) - cn^{\epsilon/2}} \rightarrow 0.$$

- A_i - the event that $\hat{\theta}^{g_i} \neq \theta_i$.

- $D(P_{\hat{\theta}_i} \| P_{\theta_i}) \geq \frac{n^{1-\epsilon}}{c}$ for $\hat{\theta} \in A_i$, c is constant.

I.I.D. Sources - Large Alphabet k - Minimax

[Shamir, 2003]

Problems with Large k

- Volume of Λ_k is $1/(k-1)!$ (decreases with n), because

$$\sum_{i=1}^{k-1} \theta_i \leq 1.$$

- Too large spacing in grid $n^{-0.5(1-\epsilon)}$ results in loose bound.

- Too small spacing $(nk)^{-0.5(1-\epsilon)}$ results in lack of distinguishability in grids.

Solution

- Build non-uniform grids.

- Spacing near $\frac{n}{\sqrt{a}}$ proportional to $\frac{n^{1-\epsilon/2}}{\sqrt{a}}$.

- Number of grid points preceding $\frac{n}{\sqrt{a}}$ proportional to $\frac{n^{\epsilon/2}}{\sqrt{a}}$.

Drawback

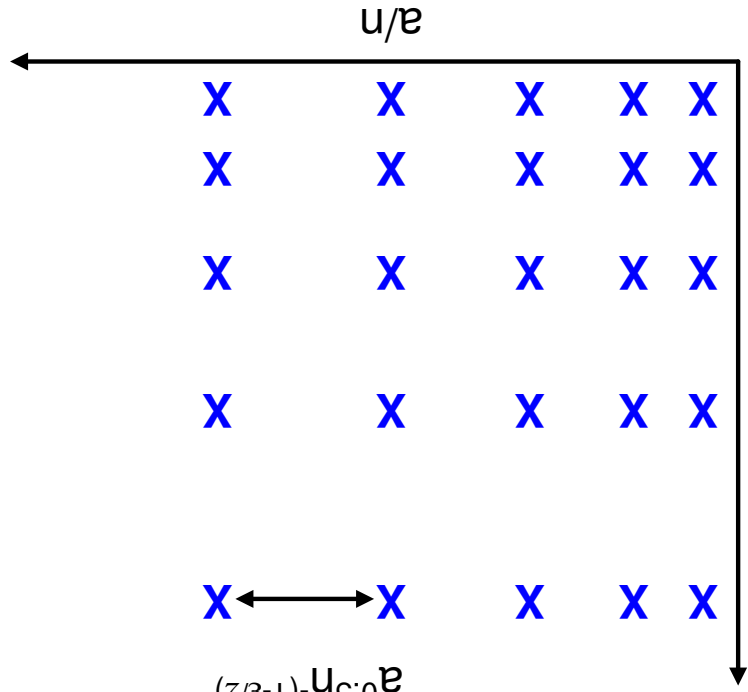
- This structure violates the requirements of the strong version, and thus is only good for minimax/maximin redundancies.

Minimax/Maximin Redundancy - I.I.D. Large k

- φ is grid below,
- $\theta \in \varphi$ distinguishable by above definition (proved as in finite parametric case),
- bounding number of points in grid results in

$$R_+^n(\Lambda_k) = R_-^n(\Lambda_k) \geq (1 - \varepsilon) \frac{2n}{(k-1)} \log \frac{k}{n}$$

$$a^{0.5n^{-(1-\varepsilon/2)}}$$



Most Sources - I.I.D. Large k

Key Realizations

- Non-uniform grid above is not useful here.
- All sources outside a $k - 1$ dimensional sphere with radius $r = n^{-0.5(1-\varepsilon)}$ around θ are distinguishable from θ by x^n .

Method

- Pack as many as possible spheres with radius r and volume $V_{k-1}(r)$ in the $k - 1$ dimensional space Λ_k of volume $1/(k - 1)!$.

- Place $\theta \in \varphi$ at centers of the spheres (whole grid shifted for random selection).
- Factor in packing density $2^{-(k-1)}$ to reduce number of points.

$$M \geq \frac{1}{V_{k-1}(r) 2^{(k-1)}}$$

Result

$R_n(L, \theta) \geq (1 - \varepsilon) \frac{2^n}{n \log k}$ for every code $L(\cdot)$ and almost every $\theta \in \Lambda_k$. [Shamir, 2003]

Note: Second order term is lower than that of min/max/maximin bound.

Patterns Induced by I.I.D. Sources

Motivation

- Classical compression considers known small alphabets.
- Sometimes alphabet is unknown and possibly large.
- Coding cost of unknown alphabet is inevitable.

Approach

- Use the inevitable cost to improve compression.
- Code sequence *patterns* in a second stage.

Patterns

- Indices assigned to original sequence letters in order of first occurrence.
- **Example:** The strings: $x_n = \text{'lossless', 'sellsoil', '12331433', '76887288'}$ all have the same pattern $\Psi(x_n) = \text{'12331433'}$.
- Individual sequence redundancy studied in [Aberg, *et al.*, 1997, Orlitsky *et al.*, 2002-].

I.I.D. Induced Patterns - Derivation

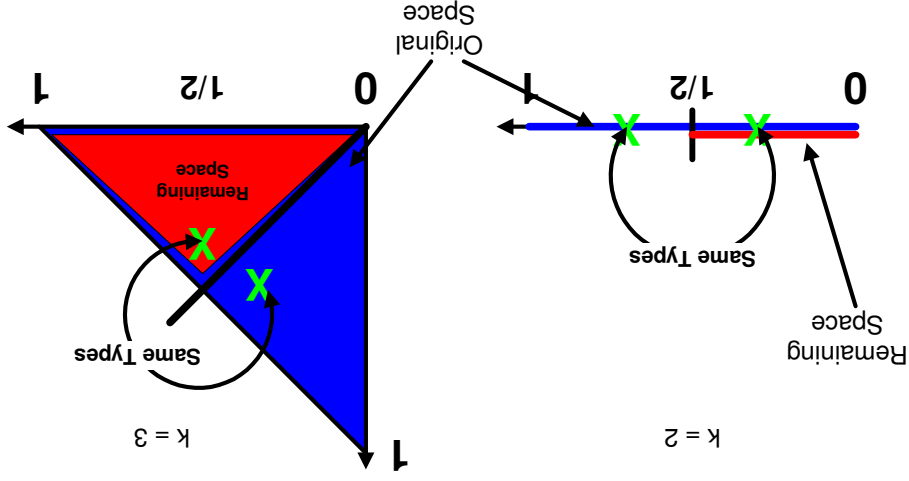
- Any θ' which is a permutation of θ appears to be the same source. Example: typical sequences - similar patterns

$$\theta = \{0.1, 0.2\} \qquad \theta' = \{0.7, 0.2\}$$

$$x_n = 1223333333 \qquad x_n = 3221111111$$

$$\Psi(x_n) = 1223333333 \qquad \Psi(x_n) = 1223333333$$

- There are at most $k!$ such permutations.



Note: for $k = 3$ this is true for any combination of 2 out of 3 letters.

Pattern Redundancy Bounds

- The grid (in both maximin and most source cases) reduces

$$M_{\Psi} \geq \frac{M_{\text{i.i.d.}}}{k!}$$

- For $k \geq n^{1/3}$ too many permutations eliminated more than once, but worst smaller k can be assumed.

- More sequences contribute to correct decision in the grid to allow

distinguishability.

Bounds [Shamir, 2003]

- Average minimax lower bound

$$R_+^n[\Psi(A^k)] \geq \left\{ \begin{array}{l} \left(\frac{2}{\pi} \right)_{1/3} \cdot (1.5 \log e) \cdot n^{-(2+\epsilon)/3} - O\left(\frac{n}{\log n}\right), \text{ for } k > \left(\frac{2}{\pi n^{1-\epsilon}}\right)_{1/3} \\ \log_{n^{1-\epsilon}}^{k_3} + \frac{2n}{k-1} \log_{\pi^{\epsilon/3}} \frac{2}{n} - O\left(\frac{n}{\log k}\right), \text{ for } k \leq \left(\frac{2}{\pi n^{1-\epsilon}}\right)_{1/3} \end{array} \right.$$

- Average most-sources lower bound

$$R_n[L, \psi(\theta)] \geq \left\{ \begin{array}{l} \log_{n^{1-\epsilon}}^{k_3} - \frac{2n}{k-1} \log_{\frac{8\pi^{\epsilon/3}}{8\pi^{\epsilon/3}}} \frac{2n}{k-1} - O\left(\frac{n}{\log k}\right), \text{ for } k \leq \frac{2}{1} \cdot \left(\frac{\pi}{n^{1-\epsilon}}\right)_{1/3} \\ \frac{1.5 \log e}{2\pi^{1/3}} \cdot n^{-(2+\epsilon)/3} - O\left(\frac{n}{\log n}\right), \text{ for } k > \frac{2}{1} \cdot \left(\frac{\pi}{n^{1-\epsilon}}\right)_{1/3} \end{array} \right.$$

Piecewise Stationary Sources - PSS's

Definition of PSS $\psi_{\triangleleft}(\theta, t) \in \Lambda^q \subset \Lambda$

- PSS - emits data divided into independent stationary segments separated by abrupt changes in statistics

- Λ - n th order class of PSS's (contains all possible combinations of the k -dimensional parameters for n -sequences)

- Λ^q - All PSS's in Λ with q segments

- $\theta_{\triangleleft} \equiv \{\theta_1, \theta_2, \dots, \theta^q\}$ - *segmental parameters*
- $t_{\triangleleft} \equiv \{t_1, t_2, \dots, t^{q-1}\}$ - *transition path* (TP)

Redundancy bound [Shamir, 2000]

$$R_n(L, \psi) \geq (1 - \varepsilon) \left(1 - b + b \frac{2}{kq + b - 1} \right) \frac{u}{\log(n/b)}$$

for every $L(\cdot)$, for almost every $\psi \in \Lambda^q$, for every q .
 in the min/max/maximin senses.

Bound Derivation - PSS's

Finite Number of Segments q

1. Λ_ϵ contains all ψ for which all segments long (longer than $n^{1-\epsilon/2}$), and all transitions are large.

2. Λ_ϵ is most class for fixed q .

3. Partition Λ_ϵ into sets as follows:

- Parse n -tuple to **phrases** of length $l = n^{1-\epsilon}$.

- For all $\psi \in \varphi$, $\forall i, t_i$ is a point in the same phrase in a grid with spacing l^ϵ .

- θ_i is a point in a grid as defined for stationary sources.

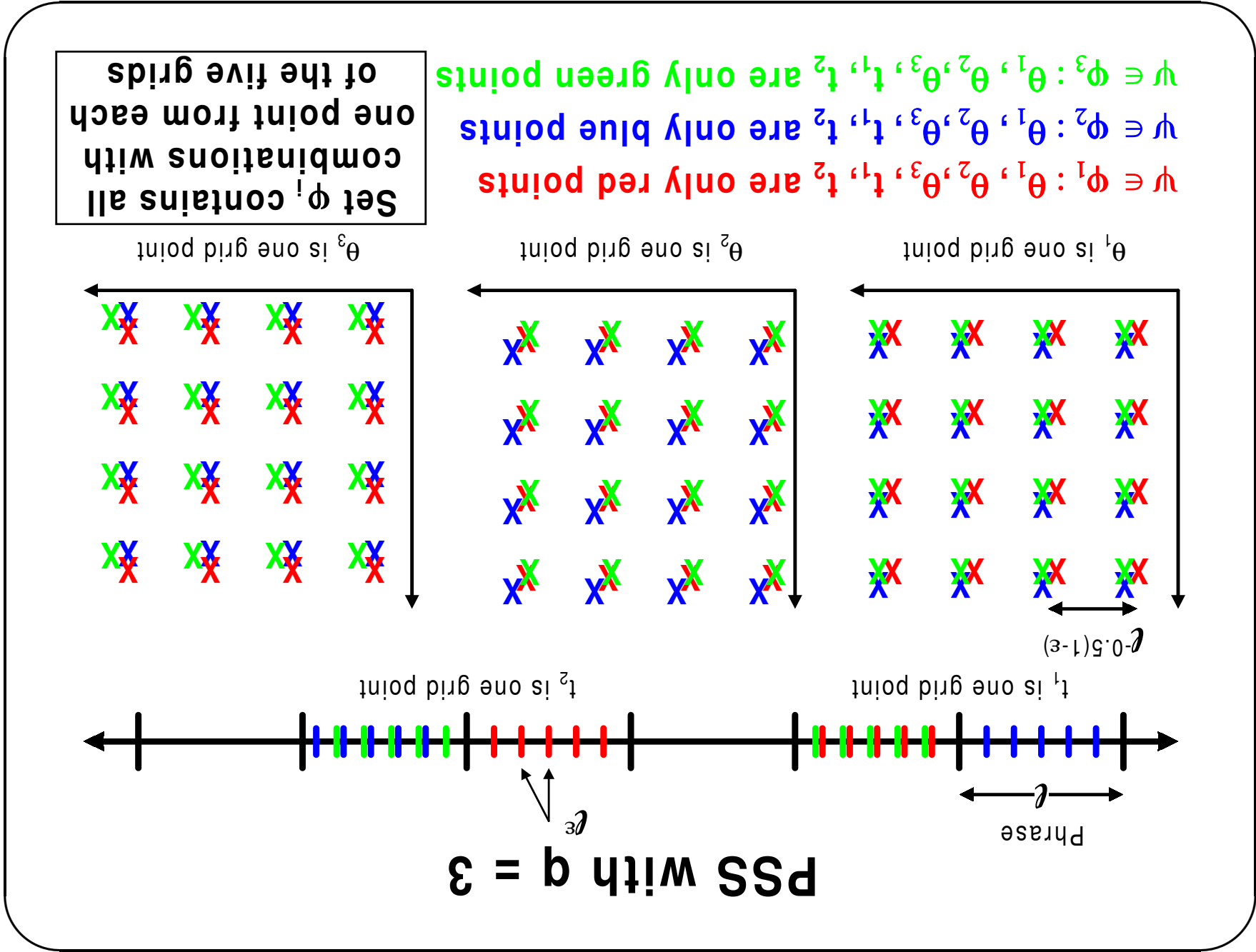
- $\forall \psi \in \varphi, t_i$ and θ_i must be from grids with identical initial shifts.

4. Distinguish among $\psi \in \varphi$ as follows:

- Use phrases entirely inside segments to estimate θ_i .

- Given $\hat{\theta}$, estimate transitions from respective grids.

By definition of the grids, the bound for finite q [Merhav, 1993] results.



General Bound Derivation - PSS's

Large q

1. Λ^ε defined is **not** most class.
2. For very large q , probability of error in **at least one** of the source parameters significantly increases the overall error probability.

Solutions to Asymptotic Problems

1. Λ^ε contains sources for which **most** segments are long and **most** transitions are large.
2. Reduce sets φ to improve distinguishability for very large q .

Two different Cases

- $q \not\gg n/q$ - almost similar to fixed q (modified according to modification 1 above).
- $q \gg n/q$ - requires additional **algebraic coding** techniques for distinguishability.

General Bound Derivation - PSS's, Cont.

Second Case: $q \gg n/q$

- Too many parameters.

- Error in estimating one results in error in estimating ψ .

Solution - Reduce φ by Linear Block Codes:

- Let $\eta > 0$ be arbitrarily small,
 - q' - number of 'free' segmental parameters,
 - c' - number of 'free' transition times.
 - $(1 - \eta) q'$ segmental parameters and $(1 - \eta) c'$ transitions chosen from grids.
 - Remaining parameters are parity checks.
 - Grids' resolutions chosen to yield Galois Fields.
 - Each grid point is assigned an element in the proper Galois Field.
 - Codes designed to correct up to $\alpha n q'$ errors (exist: Gilbert-Varshamov).
- Guarantees distinguishability even for $q \gg n/q$, resulting in the same asymptotic bound (ϵ is now larger).

Additional Source Classes

PSS's with Slowly Varying Statistics [Shamir, 2001]

- q segments, transition duration of $(n/q)^\alpha$:

$$R_n(L, \psi) \geq (1 - \varepsilon) \left[kq \frac{2}{\alpha} + (q - 1) \left(1 - \frac{2}{\alpha} \right) \right] \frac{n}{\log(n/q)}$$

- If durations unknown,

$$R_n(L, \psi) \geq (1 - \varepsilon) \left(1 - \frac{2}{\alpha} kq + q - 1 \right) \frac{n}{\log(n/q)}$$

Hierarchical version of redundancy-capacity for compound class must be used. Insignificant cost above PSS's.

Switching Sources - s states [Shamir, 2001]

If $s \leq (n/q)^{0.5k(1-\varepsilon)}$. Then, for every code $L(\cdot)$ and almost all sources

$$R_n(L, \psi) \geq (1 - \varepsilon) \frac{n}{\log(n/q)} \left[ks \frac{2}{\alpha} \log(n/s) + (q - 1) \log(n/q) + (s - b) \log s \right]$$

Otherwise,

$$R_n(L, \psi) \geq (1 - \varepsilon) \frac{n}{\log(n/q)} \left(1 - b + b \frac{2}{\alpha} kq \right)$$

Summary and Conclusions

1. The redundancy-capacity theorem is very useful to derive lower bounds on
 - min/max/redundancy in universal coding,
 - redundancy for most sources in universal coding.
2. Lower bounds on redundancy in both cases were obtained for
 - finite number of parameters parametric sources,
 - i.i.d. sources with large alphabets,
 - patterns induced by i.i.d. sources,
 - piecewise stationary sources,
 - piecewise stationary sources with slowly varying statistics,
 - switching sources.
3. Different techniques from coding theory were used:
 - random coding,
 - sphere packing,
 - algebraic code distance bounds.